

Hidden Markov Models



Dinucleotide Frequency

- Consider all 2-mers in a sequence
{AA,AC,AG,AT,CA,CC,CG,CT,GA,GC,GG,GT,TA,TC,TG,TT}
- Given 4 nucleotides: each with a probability of occurrence of $\approx \frac{1}{4}$.
Thus, one would expect that the probability of occurrence of any given dinucleotide is $\approx \frac{1}{16}$.
- However, the frequencies of dinucleotides in DNA sequences vary widely.
- In particular, CG is typically underrepresented
(frequency of CG is typically $< \frac{1}{16}$)

Example

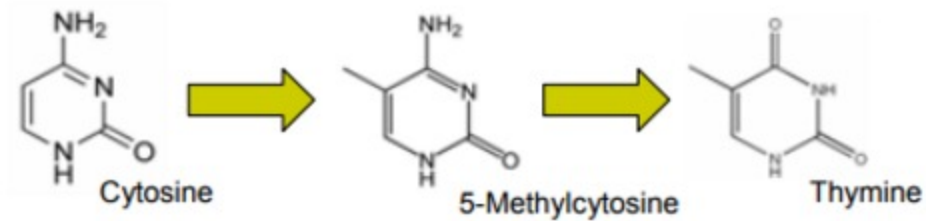
- From a 291829 base sequence

AA	0.120214646984	GA	0.056108392614
AC	0.055409350713	GC	0.037792809463
AG	0.068848773935	GG	0.043357731266
AT	0.083425853585	GT	0.046828954041
CA	0.074369148950	TA	0.077206436668
CC	0.044927148868	TC	0.056207766218
CG	0.008179475581	TG	0.063698479926
CT	0.066857875186	TT	0.096567155996

- Expected value 0.0625
- CG is 7 times smaller than expected

Why so few CGs?

- CG is the least frequent dinucleotide because C in CG is easily *methyated*. And, methylated Cs are easily mutated into Ts.



- However, methylation is suppressed around genes and transcription factor binding sites
- So, CG appears at relatively *higher* frequency in these important areas
- These localized areas of higher CG frequency are called ***CG-islands***
- Finding the CG islands within a genome is among the most reliable gene finding approaches

CG Island Analogy

- The CG islands problem can be modeled by a toy problem named "*The Fair Bet Casino*"
- The outcome of the game is determined by coin flips with two possible outcomes: Heads or Tails
- However, there are two different coins
 - A **Fair** coin:
Heads and Tails with same probability $\frac{1}{2}$.
 - The **Biased** coin:
Heads with prob. $\frac{3}{4}$,
Tails with prob. $\frac{1}{4}$.



The "Fair Bet Casino"

- Thus, we define the probabilities:
 - $P(H|Fair) = \frac{1}{2}$, $P(T|Fair) = \frac{1}{2}$
 - $P(H|Bias) = \frac{3}{4}$, $P(T|Bias) = \frac{1}{4}$
- The house doesn't want to get caught switching between coins, so they do so infrequently
- Changes between Fair and Biased coins with probability 10%



Fair Bet Casino Problem

- **Input:** A sequence $x = x_1, x_2, x_3, \dots, x_n$ of observed coin tosses made by some combination of the two possible coins (**F** or **B**).
- **Output:** A sequence $\pi = \pi_1, \pi_2, \pi_3, \dots, \pi_n$, with each π_i being either **F** or **B** indicating that x_i is the result of tossing the Fair or Biased coin respectively.



Problem Subtleties

- Any observed outcome of coin tosses *could* have been generated by *any combination* of coin exchanges
- However, all coin-exchange combinations are not equally likely.

Tosses:	T H H H H	
Coins1:	F F F F F	$P(\text{Tosses} \text{Coins1}) = 1/2 \cdot 1/2 \cdot 1/2 \cdot 1/2 \cdot 1/2 = 1/32$
Coins2:	B B B B B	$P(\text{Tosses} \text{Coins2}) = 1/4 \cdot 3/4 \cdot 3/4 \cdot 3/4 \cdot 3/4 = 81/1024$
Coins3:	F F B B B	$P(\text{Tosses} \text{Coins3}) = 1/2 \cdot 1/2 \cdot 3/4 \cdot 3/4 \cdot 3/4 = 27/256$

- We ask, "What ***coin-exchange combination*** has the highest probability of generating the observed series of tosses?"
- The coin tosses are a signal, and figuring out the most likely coin-exchange sequence is a *Decoding Problem*

Let's consider the extreme cases

- Suppose that the dealer *never* exchanges coins.
- Some definitions:
 - $P(x/Fair)$: prob. of generating the x using the Fair coin.
 - $P(x/Biased)$: prob. of generating x using the Biased coin .

$P(x|\text{fair coin})$ vs. $P(x|\text{biased coin})$

$$P(x|Fair) = P(x_1 \dots x_n|Fair) = \prod_{i=1,n} p(x_i|Fair) = \left(\frac{1}{2}\right)^n$$

$$P(x|Biased) = P(x_1 \dots x_n|Biased) = \prod_{i=1,n} p(x_i|Biased) = \left(\frac{3}{4}\right)^k \left(\frac{3}{4}\right)^{n-k} = \frac{3^k}{4^n}$$

- Where k is the number of **H**eads observed in x

$P(x|\text{Fair coin}) = P(x|\text{Biased coin})$

- When is a sequence equally likely to have come from the Fair or Biased coin?

$$P(x|\text{Fair}) = P(x|\text{Biased})$$

$$\left(\frac{1}{2}\right)^n = \frac{3^k}{4^n}$$

$$2^n = 3^k$$

$$n = k \log_2 3$$

- when $k = \frac{n}{\log_2 3}$ ($k \approx 0.63n$)
- So when the number of heads over a contiguous sequence of tosses is greater than 63% the dealer is most likely used the biased coin

Log-odds Ratio

- We can define the log-odds ratio as follows:

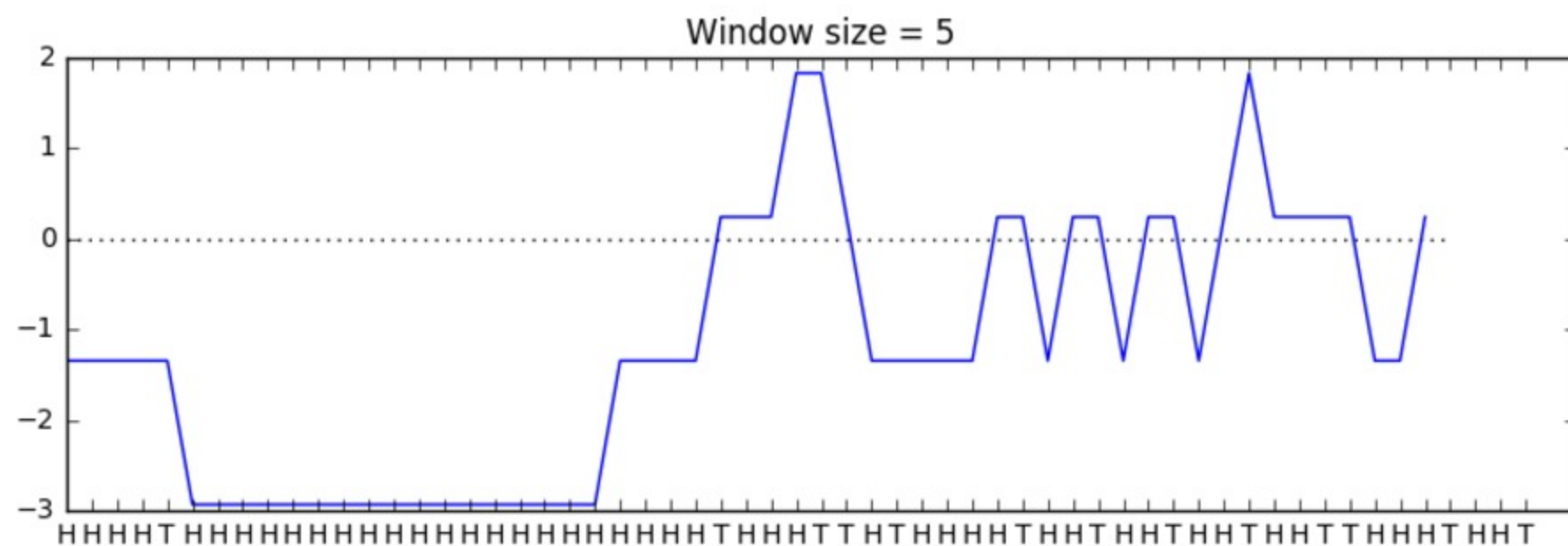
$$\begin{aligned}\log_2 \left(\frac{P(x|Fair)}{P(x|Biased)} \right) &= \sum_{i=1}^k \log_2 \left(\frac{P(x_i|Fair)}{P(x_i|Biased)} \right) \\ &= n - k \log_2 3\end{aligned}$$

- The log-odds ratio is a means (threshold) for deciding which of two alternative hypotheses is most likely
- "*Zero-crossing*" measure:
 - If the log-odds ratio is > 0 then the numerator (Fair coin) is more likely
 - if the log-odds ratio is < 0 then the denominator (Biased coin) is more likely
 - They are equally likely if the log-odds ratio = 0

Log-odds over a sliding window

- Given a sequence of length n , consider the log-odds ratio of a sliding window of length $w \ll n$

$x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, \dots, x_n$



Disadvantages:

- The length of CG-island (appropriate window size) is not known in advance
- Different window sizes may classify the same position differently
- What about the rule that they don't swap out the coins frequently?

Key Elements of the Problem

- There is an unknown or *hidden* state for each observation (Was the coin the Fair or Biased?)
- Outcomes are modeled probabilistically:
 - $P(H|Fair) = \frac{1}{2}$, $P(T|Fair) = \frac{1}{2}$
 - $P(H|Bias) = \frac{3}{4}$, $P(T|Bias) = \frac{1}{4}$
- Transitions between states are modeled probabilistically:
 - $P(\pi_i = Bias | \pi_{i-1} = Bias) = a_{BB} = 0.9$
 - $P(\pi_i = Bias | \pi_{i-1} = Fair) = a_{FB} = 0.1$
 - $P(\pi_i = Fair | \pi_{i-1} = Bias) = a_{BF} = 0.1$
 - $P(\pi_i = Fair | \pi_{i-1} = Fair) = a_{FF} = 0.9$

Hidden Markov Model (HMM)

- A generalization of this class of problem
- Can be viewed as an abstract machine with k *hidden* states that emits symbols from an alphabet Σ .
- Each state emits outputs with its own probability distribution, and the machine switches between states according to some other probability distribution.
- While in a certain state, the machine makes 2 decisions:
 - What symbol from the alphabet Σ should I emit?
 - What state should I move to next?

Why "Hidden"?

- Observers see the emitted symbols of an HMM but cannot see which state the HMM is currently in.
- Thus, the goal is to infer the most likely hidden states of an HMM based on the given sequence of emitted symbols.

HHHTHTHHTTTTHTHTHTHHHTHTHTHT
BBBFFFFFFFFFFFFFFFFBBBFFFFFFFF?



HMM Parameters

- Σ : set of emission characters.
Example: $\Sigma = \{0, 1\}$ for coin tossing
 - (0 for Tails and 1 Heads)
 - $\Sigma = \{1, 2, 3, 4, 5, 6\}$ for dice tossing
- Q : set of hidden states, emitting symbols from Σ .
 $Q = \{\text{Fair}, \text{Bias}\}$ for coin tossing

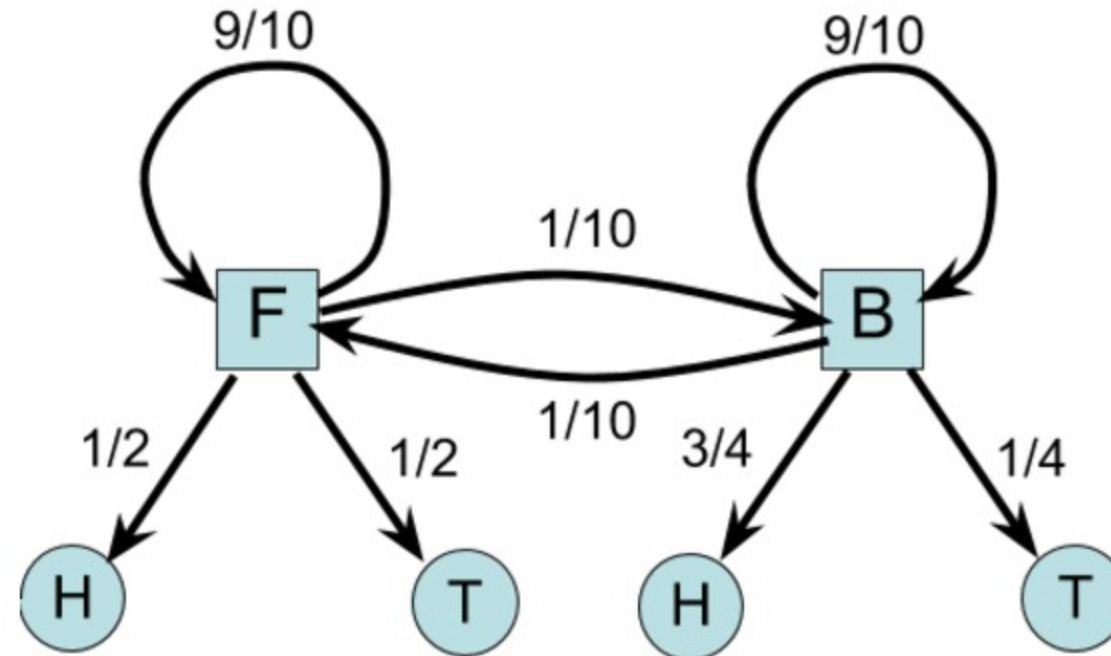
HMM for Fair Bet Casino

- The Fair Bet Casino in HMM terms:
 - $\Sigma = \{0, 1\}$ (0 for Tails and 1 Heads)
 - $Q = \{F, B\}$ – F for Fair & B for Biased coin
- Transition Probabilities A , Emission Probabilities E

A	Fair	Biased
Fair	9/10	1/10
Biased	1/10	9/10

E	Tails(0)	Heads(1)
Fair	1/2	1/2
Biased	1/4	3/4

HMM as a Graphical Model



- Directed graph with two types nodes and two types of edges
 - hidden states are shown as squares
 - emission outputs are shown as circles
 - transition edges
 - emission edges

Hidden Paths

- A path $\pi = \pi_1 \dots \pi_n$ in the HMM is defined as a sequence of hidden states.
- Consider
 - path $\pi = \text{FFFBBBBBFFF}$
 - sequence $x = \text{01011101001}$

x	=	0	1	0	1	1	1	0	1	0	0	1
π	=	<i>F</i>	<i>F</i>	<i>F</i>	<i>B</i>	<i>B</i>	<i>B</i>	<i>B</i>	<i>B</i>	<i>F</i>	<i>F</i>	<i>F</i>
$P(x_i \pi_i)$		$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{3}{4}$	$\frac{3}{4}$	$\frac{3}{4}$	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$
$P(\pi_i \rightarrow \pi_{i+1})$		$\frac{9}{10}$	$\frac{9}{10}$	$\frac{1}{10}$	$\frac{9}{10}$	$\frac{9}{10}$	$\frac{9}{10}$	$\frac{9}{10}$	$\frac{1}{10}$	$\frac{9}{10}$	$\frac{9}{10}$	

- What is the probability of the given path (FFFBBBBBFFF)?

$P(\mathbf{x}|\pi)$ calculation

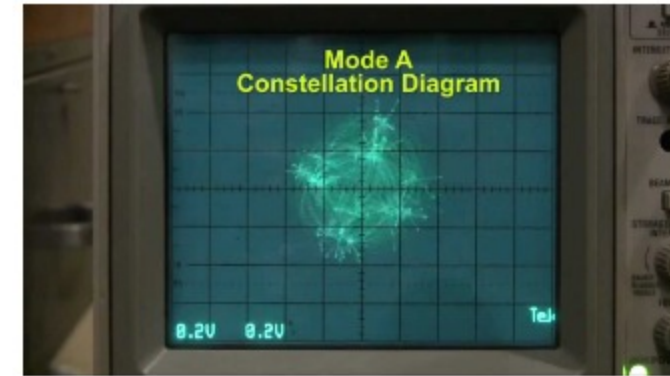
- $P(\mathbf{x}|\pi)$: Probability that sequence x was generated by the path π .

$$\begin{aligned} P(x|\pi) &= \prod_{i=1}^n P(x_i|\pi_i) \cdot P(\pi_i \rightarrow \pi_{i+1}) \\ &= \prod_{i=1}^n E_{\pi_i, x_i} \cdot A_{\pi_i, \pi_{i+1}} \end{aligned}$$

- How many such paths exist? 2^n
- What algorithmic approach would you use to find the best path? Branch and Bound? Divide and Conquer? Dynamic Programming?

Decoding Problem

Finding the optimal path in a graph is equivalent to a classical problem of decoding a message using *constellations*. This is very commonly used when a discrete set of symbols is encoded for transport over an analog medium (e.x. modems, wired internet, wireless internet, digital television).



A simple binary coding does not make good use of the dynamic range of a digital signal, however, if you put the codes too close noise becomes a problem.

- **Goal:** Find an optimal hidden path of state transitions given a set of observations.
- **Input:** Sequence of observations $x = x_1 \dots x_n$ generated by an HMM $M(\Sigma, Q, A, E)$
- **Output:** A path that maximizes $P(x|\pi)$ over all possible paths π .

How do we solve this?

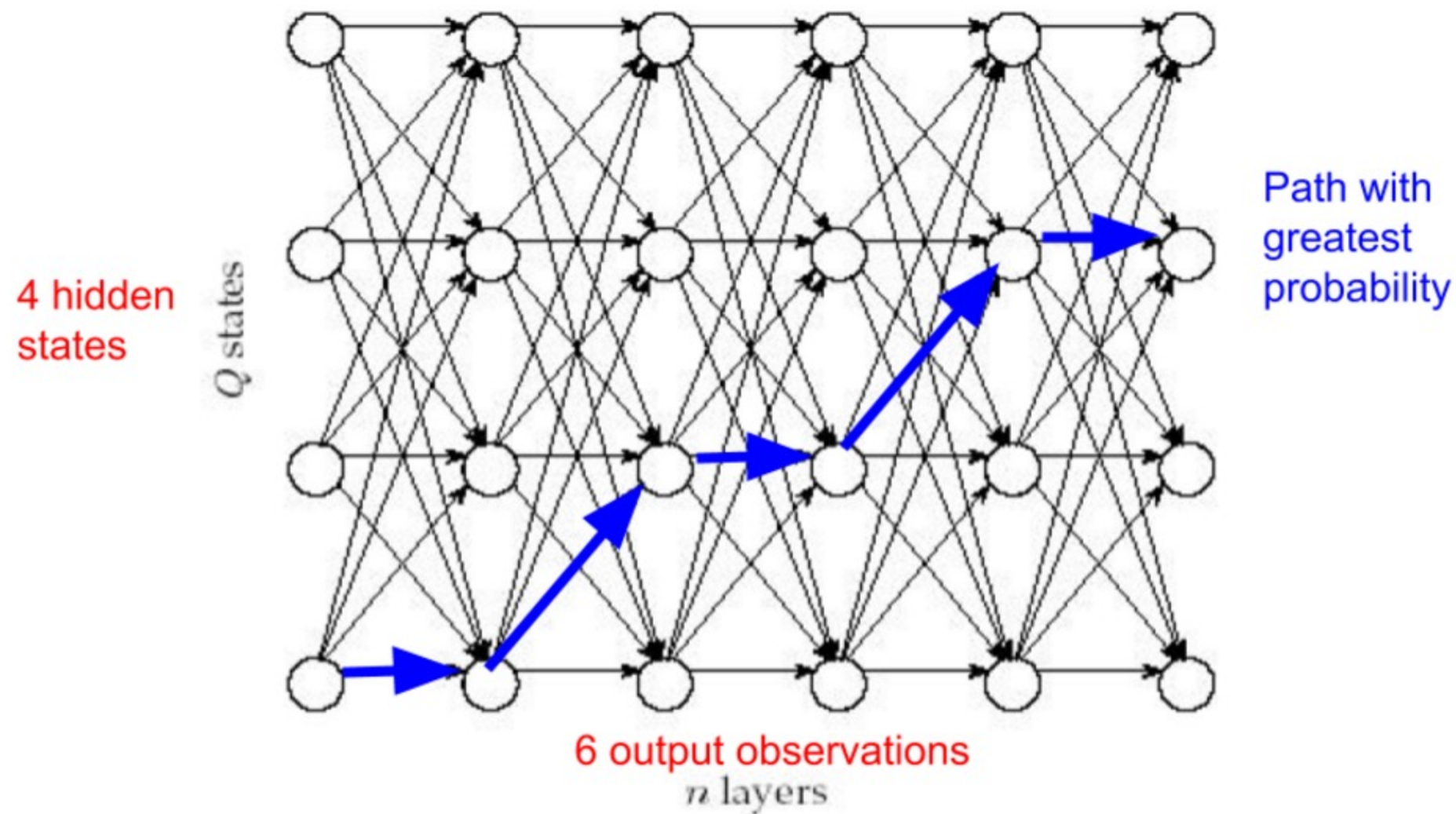
- Brute Force approach:
 - Enumerate every possible path
 - Compute $P(x_{1..n}|\pi_{1..n})$ for each one
 - Keep track of the most probable path
- A better approach:
 - Break any path in two parts, $P(x_{1..i}|\pi_{1..i}), P(x_{i..n}|\pi_{i..n})$
 - $P(x_{1..n}|\pi_{1..n}) = P(x_{1..i}|\pi_{1..i}) \times P(x_{i..n}|\pi_{i..n})$
 - Will less than the highest $P(x_{1..i}|\pi_{1..i})$ ever improve the total probability?
 - Thus to find the maximum $P(x_{1..n}|\pi_{1..n})$ we need find the maximum of each subproblem $P(x_{1..i}|\pi_{1..i})$, for i from 1 to n
 - What algorithm design approach does this remind you of?

Building Manhattan for Decoding

- In 1967, Andrew Viterbi developed a “Manhattan-like grid” (Dynamic program) model to solve the Decoding Problem.
- Every choice of $\pi = \pi_1, \pi_2, \dots, \pi_n$ corresponds to a path in the graph.
- The only valid direction in the graph is eastward.
- This graph has $|Q|^2(n - 1)$ edges.

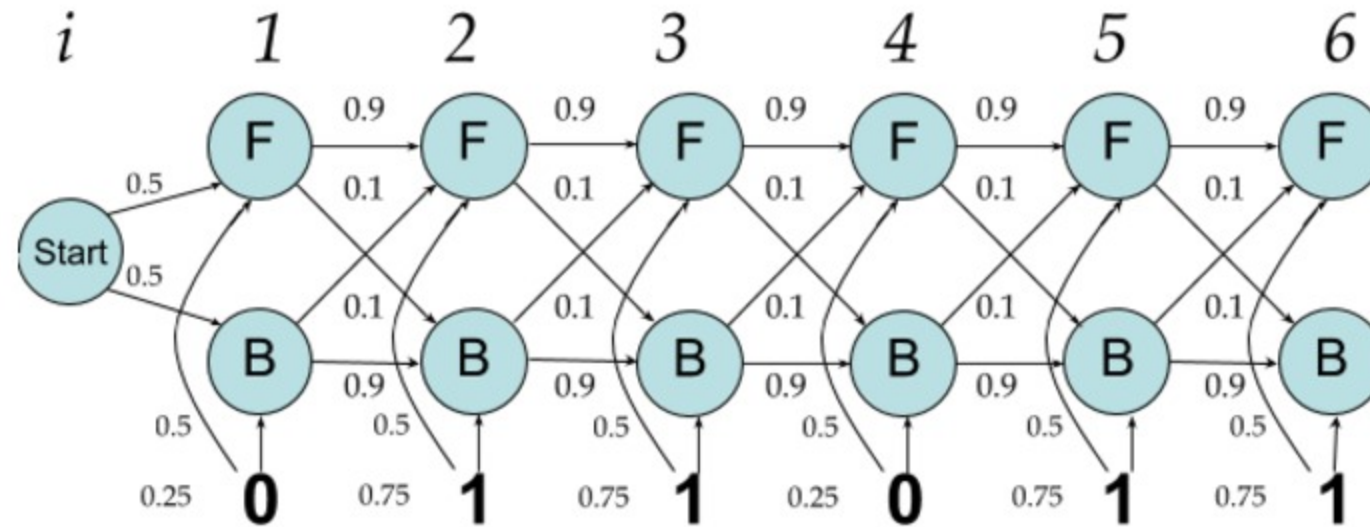


Edit Graph for Decoding Problem



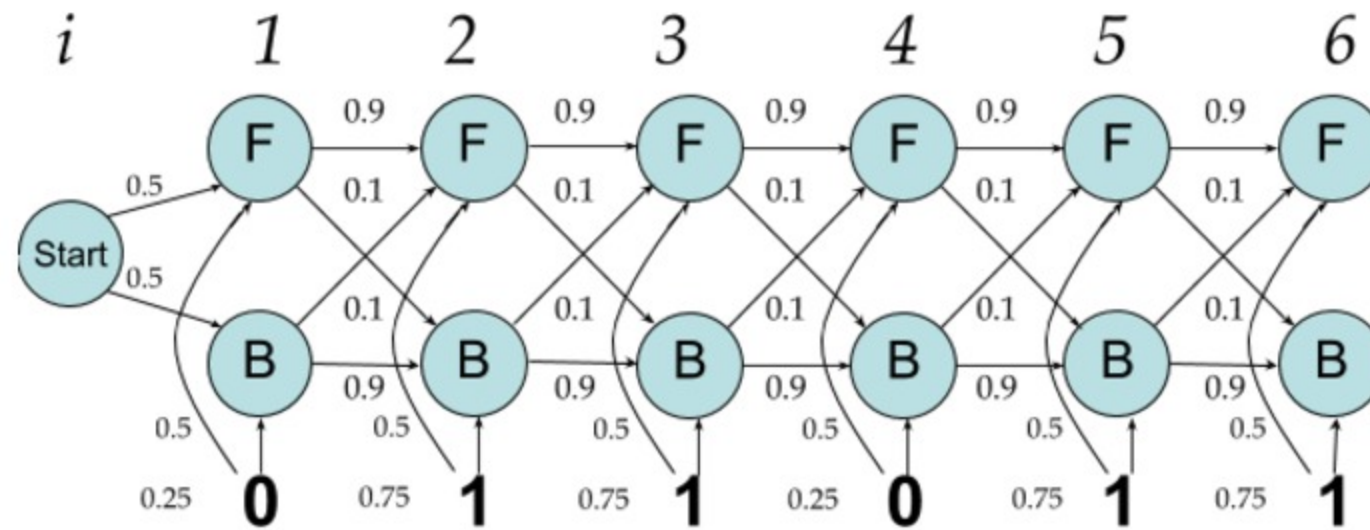
Viterbi Decoding of Fair-Bet Casino

- Each vertex represents a possible state at a given position in the output sequence
- The observed sequence conditions the likelihood of each state
- Dynamic programming reduces search space to:
 $|Q| + \text{transition_edges} \times (n-1) = 2 + 4 \times 5$ from naïve 2^6



Decoding Problem Solution

- The *Decoding Problem* is equivalent to finding a longest path in the *directed acyclic graph* (DAG), where "longest" is defined as the maximum product of the probabilities along the path.



Next Time

- We'll find the DP recurrence equations
- See examples and what it looks like in code
- See how truth and maximum likelihood do not always agree
- Apply to HMMs to problems of biological interest.